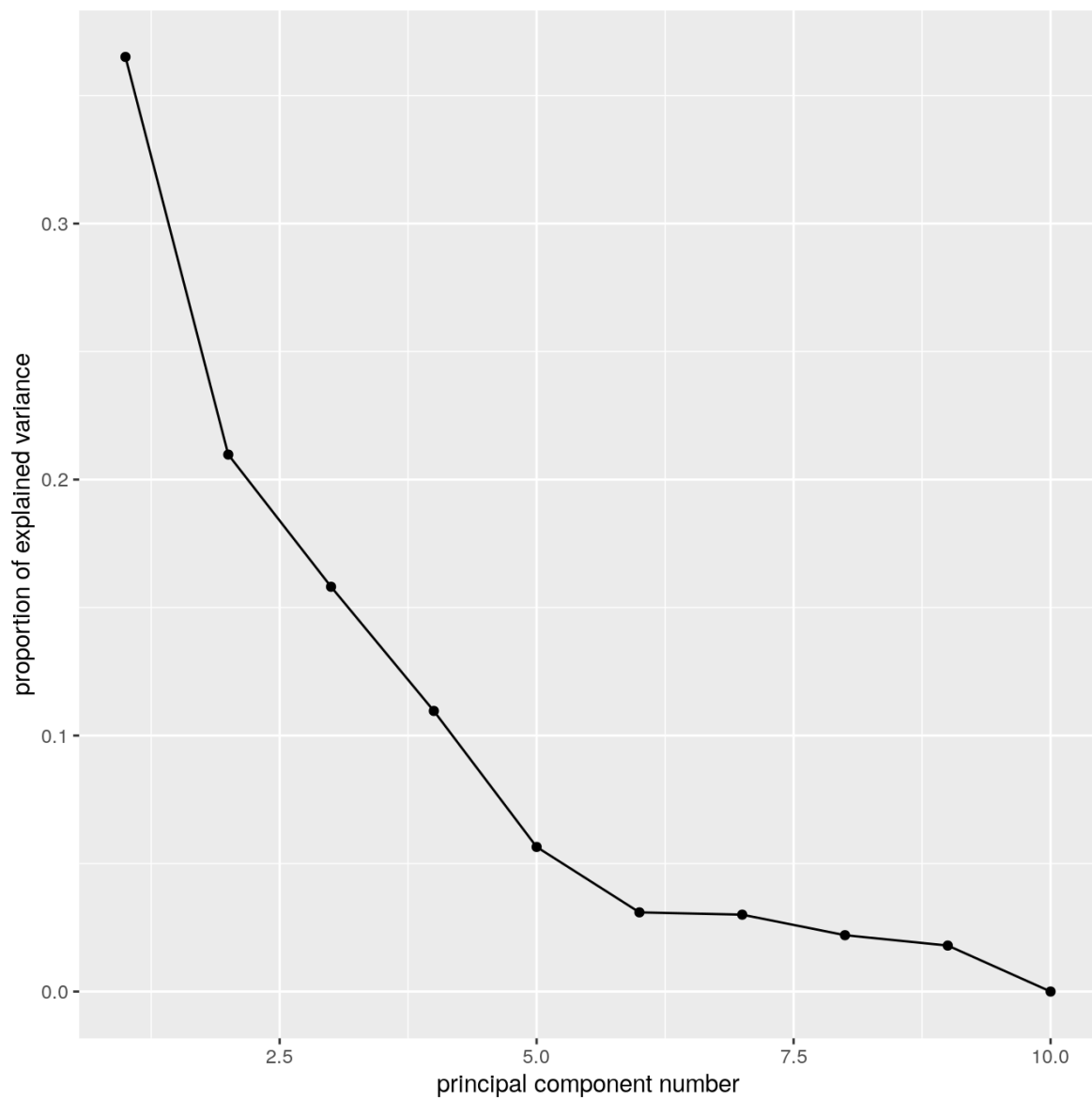# 基因组学数据分析 第三次作业

# 主成分分析法（PCA）分析流程

## 数据读入与整理

```r
library(ggplot2)
library(dplyr)
genes_expression <- read.delim(
    "./GSE106688_genes.fpkm_table.txt",
    header = TRUE
)
rownames(genes_expression) <- genes_expression[, 1]
colnames(genes_expression) <- c(
    "gene_id", "hESC R1", "hESC R2", "MES R1", "MES R2",
    "CP R1", "CP R2", "CM R1", "CM R2", "Fetal R1", "Fetal R2"
)
genes_expression <- genes_expression[, -1] # 首列作为列名
genes_expression <- genes_expression[
    rowSums(genes_expression) > 0.01,
] # 筛去全为0的行
genes_expression <- genes_expression[
    apply(genes_expression, 1, var) != 0,
] # 筛去方差为0的行
```

## PCA分析过程

```r
pca <- t(as.matrix(genes_expression)) %>%
    prcomp(
        # genes_expression,
        scale = TRUE
    ) # 此处scale参数可不加，使用单独的scale()函数可实现同样效果
var <- data.frame(pca$sdev^2)
var_per <- round(var / sum(var) * 100)
ggbiplot::ggscreeplot(pca) # 碎石图查看方差被解释的部分
```

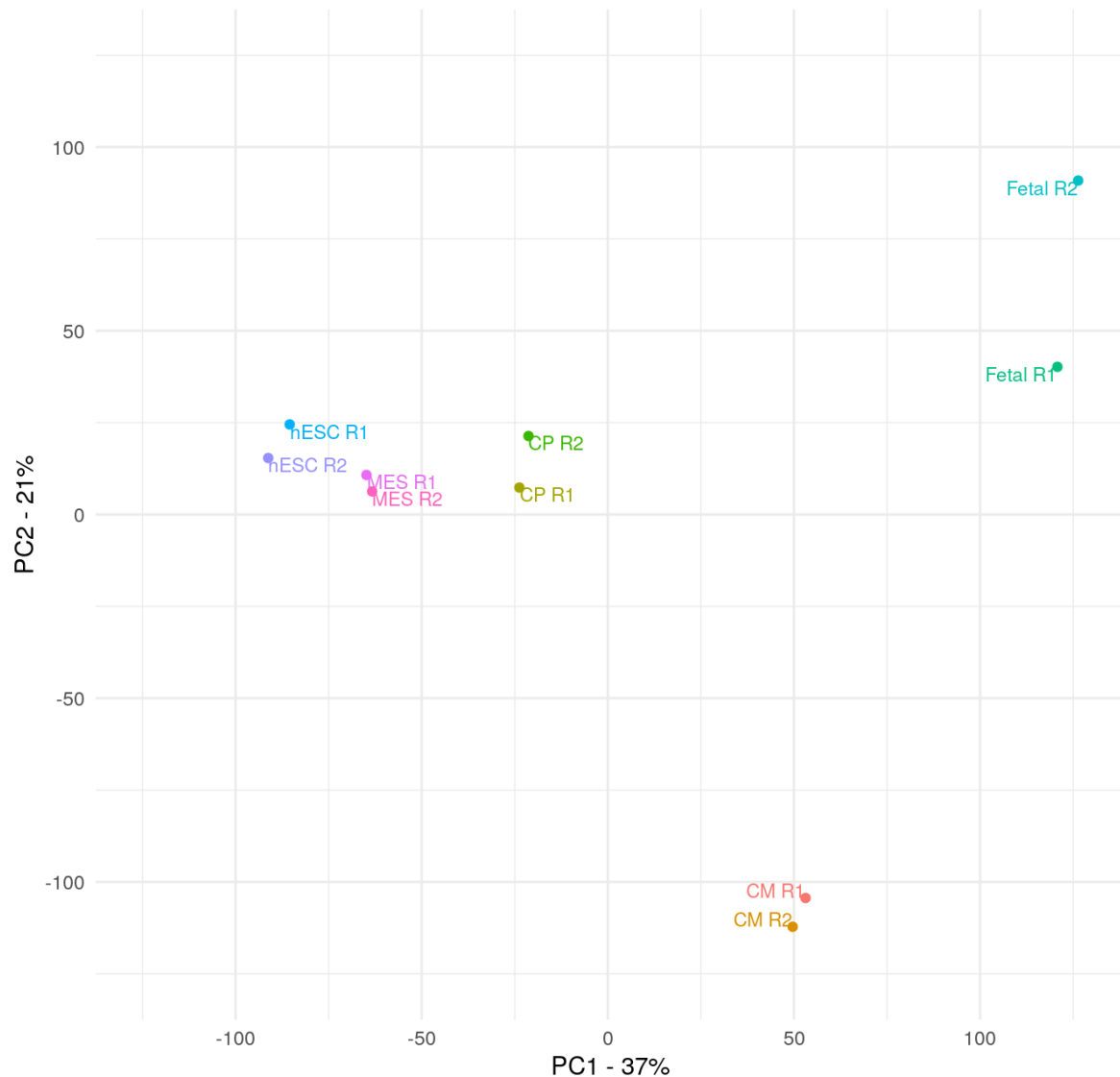## PCA结果作图

```r
data.frame(
    Sample_name = rownames(pca$x),
    X = pca$x[, 1],
    Y = pca$x[, 2]
) %>%
    ggplot(mapping = aes(
        x = X,
        y = Y,
        label = Sample_name,
        color = Sample_name
    )) +
    geom_text(
        vjust = "inward",
        hjust = "inward",
        check_overlap = TRUE,
        size = 3
    ) +
    geom_point() +
    expand_limits(x = c(-125, 125), y = c(-125, 125)) +
    theme_minimal() +
    labs(
        title = "PCA plot of the RNA-seq samples",
```

```
        x = paste("PC1 - ", var_per[1, 1], "%", sep = ""),
        y = paste("PC2 - ", var_per[2, 1], "%", sep = "")
    ) +
    theme(
        legend.position = "none"
    )
```

PCA plot of the RNA-seq samples



## PCA结果分析

### PC1

不难从图上看出，`PC1` 主要显示的是细胞类型的分化历程（

$$hESC \rightarrow MES \rightarrow CP \rightarrow CM \rightsquigarrow Fetal)$$

### PC2

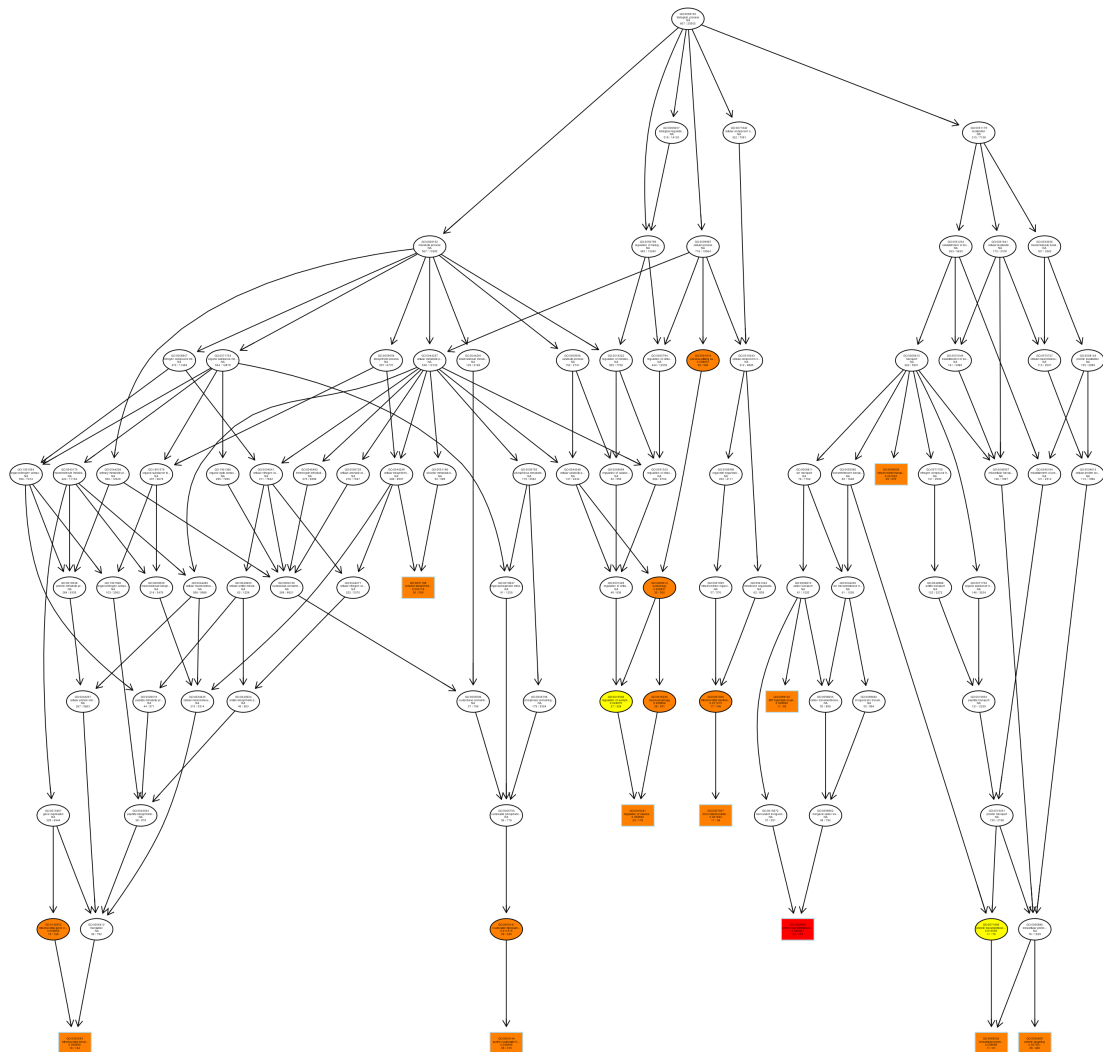`PC2` 的含义较难从PCA结果图上直观的看出，但可根据以下GO分析结果尽心推测

```
library(clusterProfiler)
library(org.Hs.eg.db)
up_bond <- unname(quantile(pca$rotation[, 2], 0.95))
low_bond <- unname(quantile(pca$rotation[, 2], 0.05))
egoP <- enrichGO(
```

```
        gene = rownames(pca$rotation[pca$rotation[, 2] > up_bond, ]),
        OrgDb = org.Hs.eg.db,
        keyType = "ENSEMBL",
        ont = "BP",
        pAdjustMethod = "BH",
        pvalueCutoff = 0.01,
        qvalueCutoff = 0.05
)
egoN <- enrichGO(
        gene = rownames(pca$rotation[pca$rotation[, 2] < low_bond, ]),
        OrgDb = org.Hs.eg.db,
        keyType = "ENSEMBL",
        ont = "BP",
        pAdjustMethod = "BH",
        pvalueCutoff = 0.01,
        qvalueCutoff = 0.05
)
plotGOgraph(egoP)
```



```
$dag
A graphNEL graph with directed edges
Number of Nodes = 51
Number of Edges = 81
```

```
$complete.dag
[1] "A graph with 51 nodes."
```

```
plotGOgraph(egoN)
```



```
$dag
A graphNEL graph with directed edges
Number of Nodes = 94
Number of Edges = 154

$complete.dag
[1] "A graph with 94 nodes."
```

**GO结果分析**

从图中不难看出，`PC2`中显著上调的基因与染色体分离、细胞核分裂、线粒体活动调控等细胞分裂的过程相关，而对应的显著下调的则是主要集中在质子转运、线粒体转录活动等过程上。因此`PC2`整体反映了细胞的分裂活动的强度，`PC2`值越大，活动强度越高，反之亦然。

# 差异表达基因热图

```r
library(pheatmap)
signif_diff_gene <- read.table("./DEGs.txt")
signif_diff_gene <- as.vector(signif_diff_gene[[1]])
diff_gene_expression <- genes_expression[rownames(genes_expression) %in%
pheatmap(
    diff_gene_expression,
    clustering_distance_rows = "correlation",
    scale = "row",
    show_rownames = FALSE,
    color = colorRampPalette(c("blue", "yellow"))(50),
    treeheight_row = 100
)
```